

Accentuation and compatibility: Replication and extensions of Shafir (1993) to rethink choosing versus rejecting paradigms

Subramanya Prasad Chandrashekar^{*†} Jasmin Weber^{‡†}

Sze Ying Chan^{‡†} Won Young Cho^{‡†} Tsz Ching Connie Chu^{‡†}

Bo Ley Cheng[‡] Gilad Feldman^{‡†§}

Abstract

We conducted a replication of Shafir (1993) who showed that people are inconsistent in their preferences when faced with choosing versus rejecting decision-making scenarios. The effect was demonstrated using an enrichment paradigm, asking subjects to choose between enriched and impoverished alternatives, with enriched alternatives having more positive and negative features than the impoverished alternative. Using eight different decision scenarios, Shafir found support for a compatibility principle: subjects chose and rejected enriched alternatives in choose and reject decision scenarios ($d = 0.32$ [0.23,0.40]), respectively, and indicated greater preference for the enriched alternative in the choice task than in the rejection task ($d = 0.38$ [0.29,0.46]). In a preregistered very close replication of the original study ($N = 1026$), we found no consistent support for the hypotheses across the eight problems: two had similar effects, two had opposite effects, and four showed no effects (overall $d = -0.01$ [-0.06,0.03]). Seeking alternative explanations, we tested an extension, and found support for the accentuation hypothesis.

Keywords: pre-registered replication; choosing versus rejecting; compatibility principle; replication crisis, accentuation hypothesis

^{*}Lee Shau Kee School of Business and Administration, The Open University of Hong Kong, Hong Kong SAR. ORCID: 0000-0002-8599-9241

[†]Contributed equally, joint first authors

[‡]Department of Psychology, University of Hong Kong, Hong Kong SAR. Email: gfeldman@hku.hk.

[§]Corresponding author. Email: gfeldman@hku.hk. ORCID: 0000-0003-2812-6599

Authorship declaration: Gilad led the reported replication effort with the team listed below. Gilad supervised each step of the project, conducted the pre-registration, and ran data collection. Prasad followed up on initial work by the other coauthors to verify and conduct additional analyses, and completed the manuscript draft. Prasad and Gilad jointly finalized the manuscript for submission. Jasmin Weber, Chan Sze Ying, Won Young Cho, and Chu Tsz Ching (Connie) conducted the replication as part of a university course.

1 Introduction

Early rational choice theories assumed the principle of invariance in human decision-making (von Neumann & Morgenstern, 1947). The invariance principle states that people's preference does not change when a decision task is described differently (description invariance) or when there are variations in the elicitation procedure (procedural invariance). Daniel Kahneman, Amos Tversky, and colleagues demonstrated that the assumptions of both description invariance and procedural invariance are often violated in human decision-making. For example, Tversky and Kahneman's (1986) findings on framing effects demonstrated that the invariance principle is violated when decision scenarios are described in a positive or negative frame. Similarly, variations in elicitation procedures were shown to cause preference reversals during the selection of job candidates (Tversky, Sattath & Slovic, 1988) and in the prediction of others' academic performance (Slovic, Griffin & Tversky, 1990).

Shafir (1993) was the first to employ the enrichment paradigm to further demonstrate the violations of procedural invariance. His study contrasted two decision-making scenarios that are intuitively equivalent: choosing versus rejecting. Subjects were randomly assigned to either choosing the preferred option from two alternatives or rejecting the option not preferred among two alternatives. The choice sets consisted of an enriched option, with both positive and negative features, and an impoverished alternative that had neutral features. Across eight scenarios, the original study found that the enriched alternative was selected more often in a choice task and was rejected more often in a rejection task. Shafir interpreted the results based on the compatibility principle that predicts that decision outcomes depend on the weighing of positive features during a choice task and negative features during a rejection task. That is, decision-makers focus their attention on positive features during a choice task as they need positive reasons to justify the choice, whereas they direct their attention to negative features during the rejection task as they need reasons to reject an alternative. We summarized the scenarios used in the original article in Table 1 and the findings in Table 2 and Table 3.

They conducted an initial analysis of the paper, designed the replication, initiated the extensions, wrote the pre-registration, conducted initial data analysis, and wrote initial replication reports. Bo Ley Cheng guided and assisted the replication effort in the course.

Acknowledgments: Subramanya Prasad Chandrashekar would like to thank the Institute of International Business and Governance (IIBG), established with the substantial support of a grant from the Research Grants Council of the Hong Kong Special Administrative Region, China (UGC/IDS 16/17), for its support.

Copyright: © 2021. The authors license this article under the terms of the Creative Commons Attribution 3.0 License.

TABLE 1: Summary of scenarios in Shafir (1993) Experiments 1 to 8.

Problem	Scenario	Impoverished alternative	Enriched alternative
1	Which parent to award/deny the sole custody of the child	Parent A	Parent A
2	Which vacation spot to prefer/cancel	Spot A	Spot B
3	Which course to take immediately/postpone	Course X	Course Y
4	Which lottery to choose/give up	Lottery 1	Lottery 2
5	Which lottery to choose/give up	Lottery 1	Lottery 2
6	Which ice cream flavor to choose/give up	Flavor A	Flavor B
7	Which candidate to vote for/not to vote for	Candidate A	Candidate B
8	Which lottery to choose/Which lottery to reject first, and then reject later	Lottery 1, Lottery 2	Lottery 3

1.1 Choice of target article: Shafir (1993)

Shafir's (1993) article has been highly influential, with more than 640 citations, and has contributed to an active literature on the relational properties of choice sets. The compatibility principle has formed the theoretical basis for explaining people's decisions when deciding between products and job applicants (Park, Jun & MacInnis, 2000; Sokolova & Krishna, 2016), and when choosing among products (Chernev, 2009; Nagpal & Krishnamurthy, 2008). Furthermore, the findings of the original article have formed the basis for subsequent theoretical work (e.g., Kahneman, 2003; Morewedge & Kahneman, 2010; Shafir, Simonson & Tversky, 1993).

Recently, Many Labs 2 (Klein et al., 2018) conducted a partial replication of Problem 1 from the original study. The findings of this partial replication failed to provide support for the compatibility principle and the original findings. In response to the replication, Shafir (2018) noted several limitations in the replication effort. First, Klein and co-authors attempted to replicate only one out of eight decision scenarios reported in the original study. Second, the nature and value of the alternatives presented in the chosen decision-making scenario may have changed in meaning since the publication of the original study due to societal changes over time. Third, unlike the original study, the replication study did not counterbalance the order of presentation of the alternatives. In the current replication, conducted before we were aware of the Many Labs effort, we addressed the noted methodological limitations of the earlier replication (Shafir, 2018), and went beyond the replication to add extensions to try and gain further insights about the phenomenon. More on that below.

TABLE 2: Descriptive statistics: The percentages of subjects who Chose/Rejected across all problems in the original study and current replication.

Problem Options	Original study			Replication ($N = 1026$)			
	N	Choose-Group	Reject-Group	Choose + Reject	Choose-Group	Reject-Group	Choose + Reject
1 Parent A (I)	170	36%	45%	81%	45.70%	52.50%	98.20%
Parent B (E)		64%	55%	119%	54.30%	47.50%	101.80%
2 Spot A (I)	172	33%	52%	85%	55.60%	48.60%	104.20%
Spot B (E)		67%	48%	115%	44.40%	51.40%	95.80%
3 Course X (I)	424	25%	65%	90%	45.50%	59.40%	104.90%
Course Y (E)		75%	35%	110%	54.50%	40.60%	95.10%
4 Lottery 1 (I)	279	25%	50%	75%	18.30%	68.80%	87.10%
Lottery 2 (E)		75%	50%	125%	81.70%	31.30%	113.00%
5 Lottery 1 (I)	278	23%	60%	83%	14.80%	66.60%	81.40%
Lottery 2 (E)		77%	40%	117%	85.20%	33.40%	118.60%
6 Flavor A (I)	359	28%	55%	83%	43.60%	53.10%	96.70%
Flavor B (E)		72%	45%	117%	56.40%	46.90%	103.30%
7 Candidate A (I)	398	79%	8%	87%	90.70%	29.10%	119.80%
Candidate B (E)		21%	92%	113%	9.30%	70.90%	80.20%
8 Lottery 1 (I)	139	39%	44%	83%	10.90%	41.30%	52.20%
Lottery 2 (I)		21.00%	36.90%	57.90%			
Lottery 3 (E)		61%	56%	117%	68.10%	21.80%	89.90%

Note: (I) = Impoverished option, (E) = Enriched option. In the replication, subjects ($N = 1026$) completed all 8 problems.

We choose to conduct a replication of Shafir (1993) due to its impact (Coles, Tiokhin, Scheel, Isager & Lakens, 2018; Isager, 2019), aiming for a comprehensive independent replication of all problems in the article. Replications are especially relevant following the recent recognition of the importance of reproducibility and replicability in psychological science (e.g., Brandt et al., 2014; Open Science Collaboration, 2015; van't Veer & Giner-Sorolla, 2016; Zwaan, Etz, Lucas & Donnellan, 2018). A comprehensive replication of this target article is needed, given the ongoing discussion regarding the evaluation of replications and the active debate around the findings of Many Labs 2 and other mass-replication efforts.

Our predictions in the replication followed that of Shafir (1993):

Hypothesis 1: Subjects choose and reject the enriched alternative more often than the impoverished alternative across task frames (choice vs. rejection).

TABLE 3: Summary of findings comparing the original study’s and the replication’s results.

Problem	Shafir (1993)	Replication					Replication summary
	Cohen’s <i>d</i>	<i>z</i>	<i>p</i>	Cohen’s <i>d</i>	BF ₁₀	BF ₀₁	
Hypothesis 1							
1	0.39 [0.08,0.69]	0.56	.287	0.04 [-0.09,0.16]	0.13	7.70	NS-I
2	0.32 [0.02,0.62]	-1.37	.915	-0.09 [-0.21,0.04]	0.03	29.38	NS-IO
3	0.21 [0.02,0.40]	-1.56	.941	-0.10 [-0.22,0.02]	0.03	31.92	NS-IO
4	0.51 [0.27,0.75]	4.18	.001	0.26 [0.14,0.39]	951.66	0.00	S-IW
5	0.34 [0.11,0.58]	5.99	.001	0.38 [0.26,0.50]	9.74×10 ⁶	0.00	S-C
6	0.34 [0.14,0.62]	1.06	.144	0.07 [-0.06,0.19]	0.23	4.29	NS-I
7	0.23 [0.04,0.43]	-6.37	.999	-0.41 [-0.53, -0.28]	0.01	104.42	NS-IO
8	0.34 [0.01,0.68]	-10.00	1.00	-0.53 [-0.63, -0.43]	0.01	197.93	NS-IO
Hypothesis 2							
1	0.39 [0.09,0.69]	0.56	.288	0.03 [-0.09,0.16]	0.21	4.85	NS-I
2	0.31 [0.01,0.61]	-1.37	.915	-0.09 [-0.21,0.04]	0.05	18.51	NS-IO
3	0.22 [0.03,0.41]	-1.58	.944	-0.10 [-0.22,0.02]	0.05	19.89	NS-IO
4	0.53 [0.30,0.77]	4.81	.001	0.30 [0.18,0.43]	2.64×10 ⁴	0.00	S-C
5	0.37 [0.13,0.61]	6.97	.001	0.45 [0.32,0.57]	9.87×10 ⁹	0.00	S-C
6	0.50 [0.29,0.71]	1.06	.144	0.07 [-0.06,0.19]	0.38	2.65	NS-I
7	0.38 [0.18,0.57]	-8.04	1.00	-0.52 [-0.64, -0.39]	0.02	64.33	NS-IO
8	0.34 [0.01,0.68]	-4.32	.999	-0.22 [-0.32, -0.12]	0.02	46.16	NS-IO

Note: Cohen’s *d* value is presented with 95% confidence interval within the brackets; BF = Bayes factor; Replication summary based on LeBel, Vanpaemel, Cheung, and Campbell (2019): NS-I= No signal – inconsistent; NS-IO= No signal – inconsistent (opposite); S-IW=Signal – inconsistent (weaker); S-C= Signal – consistent.

Hypothesis 2: Subjects prefer the enriched alternative more often in the choice task frame than during the rejection task frame.

1.2 Extension: Accentuation hypothesis

There were other findings and theoretical accounts for the choosing versus rejecting paradigm. Ganzach (1995) reported results opposite to Shafir (1993) by showing that preference for the enriched alternative was greater in the rejection than in the choice condition. Wedell (1997) proposed a theoretical resolution of the inconsistent findings by Shafir (1993) and Ganzach (1995). Wedell’s (1997) accentuation hypothesis stated that there is a greater need for justification in the choice condition than in the reject condition, and attribute differences are weighted more strongly when choosing due to a greater need for justification. If the enriched alternative was overall more attractive than the impoverished alternative, the positive differences were accentuated, and it was preferred more when choosing than when rejecting. If the enriched alternative was overall less attractive, negative differences were accentuated, and it was rejected more when choosing than when rejecting. This was noted

by Shafir (2018) as one of the "more interesting possibilities" for the replication failure in Many Labs 2 (p. 495).

We extended the original design by examining the attractiveness of each of the alternatives in a choice set. Unlike binary choice, continuous scales for each option allow for higher sensitivity (how far apart are the differences in preferences between alternatives) and advanced analyses to examine the alternative theoretical explanation. Using these measures, we were able to run analyses to test the accentuation hypothesis.

1.3 Extension: Choice ability and preferences predictor

We also added an extension to examine the association between trait choice ability and preference and choosing versus rejecting decisions. Previous research on choice has argued that choice mindset, a psychological tendency, is associated with people ascribing agency to themselves and perceiving their own and others' actions through the lens of choice (Savani, Markus, Naidu, Kumar & Berlia, 2010). People with a choice mindset view mundane actions such as checking emails and reading newspapers not as mere actions, but as choices. Thus, people with a choice mindset are prone to approach decisions with a clear choice framework. Building on the compatibility principle, we expected that individuals who rate themselves high on the ability to choose and indicate a high preference for choice would be more likely to prefer enriched alternatives, because they are more likely to take on the choosing strategy over the rejecting strategy in comparison to people who rate themselves lower on the ability to choose and indicate a low preference for choice.

2 Method

2.1 Pre-registration, power analysis, and open science

We preregistered the experiment on the Open Science Framework (OSF). Disclosures, power analyses, all materials, and additional details and analyses are available in the Supplementary Material. All measures, manipulations, and exclusions are reported, and data collection was completed before analyses. Pre-registration is available at: <https://osf.io/r4aku>. Data and R/RMarkdown code (R Core Team, 2015) are available at: <https://osf.io/ve9bg/>. We preregistered with the aim of detecting the smallest effect size ($d = 0.21$) observed in the original study at power of 0.95, which suggested a sample size of 1092.

2.2 Subjects

A total of 1026 subjects were recruited online through American Amazon Mechanical Turk (MTurk) using the TurkPrime.com platform (Litman, Robinson & Abberbock, 2017) ($M_{age} = 39.39$, $SD_{age} = 12.47$; 540 females). In the pre-registration stage, we planned to report full sample findings and to examine possible exclusion criteria such as self-reported seriousness,

English proficiency, and failing attention checks. We found that exclusions had no impact on the findings.

2.3 Procedure

After consenting to take part in the study, subjects answered measures on their general attitudes towards choice (the extension). Subjects were then randomly assigned to one of two between-subject experimental conditions, either to choose (award or indicate a preference for) an option or to reject (deny or give up) an option. Each of the two experimental conditions consisted of eight decision problems (summarized in Table 1). Seven of the eight problems presented to all the subjects included a choice between two alternatives (binary; Problems 1–7) and one problem consisted of three alternatives (non-binary; Problem 8). Problems with binary alternatives had one option with both more positive and negative aspects (enriched alternative) and one with fewer positive and negative features (impoverished alternative). The problem with non-binary alternatives included one enriched alternative and two impoverished alternatives. For the non-binary problem (Problem 8), half the subjects were asked to choose a lottery that they most preferred among three alternatives, and another half in a two-step decision rejected lotteries that they least preferred, rejecting one at a time. All descriptions and questions were taken from the original article (Shafir, 1993). A comparison of the original study's sample and the replication sample is provided in Table 4 (see Table S1, in which we note the reasons for the chosen differences between original studies and the replication attempt).

TABLE 4: Comparison between original and the replication study.

	Original study	Replication
Number of problems	8 problems that included 7 binary-problems and 1 non-binary problem	8 problems that included 7 binary-problems and 1 non-binary problem
Design	Between-subjects: Design followed two between-subjects conditions for each of the binary and non-binary problems	Between-subjects: Design followed two between-subjects conditions for each of the binary and non-binary problems
Procedure	Conducted in a lab using paper and pencil. Subset of 2-3 problems out of the set, separated by filler items.	Conducted online using Qualtrics. All 8 problems, with no filler items.
Sample size	Ranged between 139 to 424 per problem across 8 problems	1026
Sample population	Undergraduates a university in USA	Subjects from Amazon Mechanical Turk (MTurk)
Remuneration	Monetary reward	Monetary reward

2.4 Measures

Trait choice (ability and preference). Two items measured the subjects' perceived ability to choose: "It's very hard for me to choose between many alternatives." (reversed) and "When faced with an important decision, I prefer that someone else chooses for me." (reversed) ($\alpha = .63$). Similarly, subjects rated their preference toward choice in two items: "The more choices I have in life, the better" and "In each decision I face, I prefer to have as many alternatives as possible to choose from." ($\alpha = .81$) On all four items, the scale was from 1 (*Strongly Disagree*) to 7 (*Strongly Agree*). The two scales were adapted from Feldman, Baumeister and Wong (2014).

Attractiveness. For each of the eight problems, after choosing or rejecting the alternative(s), subjects proceeded to the next page and rated the relative attractiveness of the enriched and impoverished alternatives. As the term "attractive" might be associated with choosing, this may lead to biases in the ratings, thus subjects were asked to rate each alternative with the terms "bad" and "good" to maintain neutrality. The scale for the items ranged from 0 (*Very bad*) to 5 (*Very good*).

2.5 Data analysis plan

We employed one-proportion and two-proportion z-tests to investigate Hypothesis 1 and Hypothesis 2, respectively. Given the clear directionality of the predictions in the original study, both tests were one-tailed. We then used the obtained z-value to calculate the Cohen's *d* effects and 95% confidence intervals. All analyses were performed using the R programming environment (R Core Team, 2015). Furthermore, we complemented Null Hypothesis Significance Testing (NHST) analyses with Bayesian analyses to quantify support for the null hypothesis when relevant (Kruschke & Liddell, 2018; Vandekerckhove, Rouder & Kruschke, 2018) using the 'BayesFactor' R package (Version 0.9.12–4.2; Morey & Rouder, 2015) and 'abtest' R package (Version 0.1.3.; based on a model by Kass & Vaidyanathan, 1992). The Bayesian analyses were added after preregistering the data analysis plan.

2.6 Evaluation criteria for replication design and findings

Table 5 provides a classification of this replication using the criteria by LeBel, McCarthy, Earp, Elson and Vanpaemel (2018) (also see Figure S1). We summarized the current replication as a "very close replication". To interpret the replication results, we followed the framework by LeBel et al. (2019). They suggested a replication evaluation using three factors: (a) whether a signal was detected (i.e., confidence interval for the replication effect size (ES) excludes zero), (b) consistency of the replication ES with the original study's ES, and (c) precision of the replications ES estimate (see Figure S2).

TABLE 5: Classification of the two replication studies based on LeBel et al.'s (2018) taxonomy.

Design facet	Replication study
IV operationalization	Same
DV operationalization	Same
IV stimuli	Same
DV stimuli	Same
Procedural details	Similar (minor adjustments)
Physical settings	Different
Contextual variables	Different
Replication classification	Very close replication

Note: Information on this classification is provided in LeBel et al. 2018. See also figure provided in the Supplementary Material.

3 Results

The proportions of subjects choosing or rejecting the enriched or impoverished alternative in each of the eight problems are detailed in Table 2. The findings of the statistical tests and effect-size estimates are summarized in Table 3 (also see Figures 1 and 2).

The enriched alternatives share exceeded 100% for Problems 1, 4, 5 and 6 (Table 2) across both the choosing and rejecting frames. The results of one-proportion z-test investigating Hypothesis 1 indicated support in Problem 4 ($z = 4.18, p < .001, d = 0.26, 95\%$ CI [0.14,0.39]) and Problem 5 ($z = 5.99, p < .001, d = 0.38, CI [0.26,0.50]$). The results were in the opposite direction for Problem 7 ($z = -6.37, p = 1.00, d = -0.41, CI [-0.53,-0.28]$) and Problem 8 ($z = -10.00, p = 1.00, d = -0.53, CI [-0.63,-0.43]$). The results of Problem 1, 2, 3, and 6 failed to provide empirical support for the compatibility hypothesis (effect sizes ranged from -0.10 CI $[-0.22,0.02]$ to 0.07 CI $[-0.06,0.19]$). Unlike the original study, these findings do not indicate consistent evidence in support of the Hypothesis 1 prediction that the enriched alternative is selected and rejected more often.

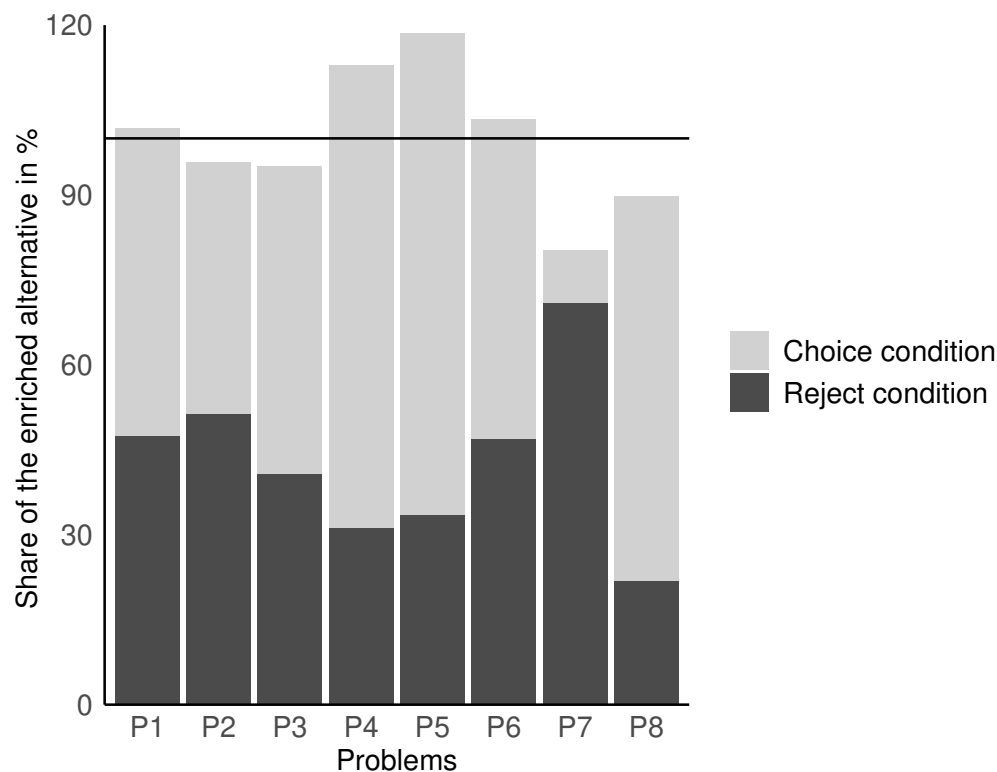


FIGURE 1: Share (in percentage) of the enriched alternative chosen and rejected across 'choice' and 'reject' experimental conditions, respectively.

We complemented the NHST analyses used in the original article with Bayesian analysis to allow for quantifying the evidence in support of the null hypothesis (see Table 3). We conducted one-sided Bayesian tests of single proportions with a prior r scale set at 0.5 (defined as "medium" and considered the more conservative option). The result revealed that Bayes factor (BF) for Problems 1, 2, 3, 6, 7, and 8 was in stronger support of the null: Problem 1: $BF_{10} = 0.13$, $BF_{01} = 7.7$; Problem 2: $BF_{10} = 0.03$, $BF_{01} = 29.38$; Problem 3: $BF_{10} = 0.03$, $BF_{01} = 31.92$; Problem 6: $BF_{10} = 0.23$, $BF_{01} = 4.29$; Problem 7: $BF_{10} = 0.01$, $BF_{01} = 104.42$; Problem 8: $BF_{10} = 0.01$, $BF_{01} = 197.93$. For example, the Bayes factor (BF_{01}) of 7.7 in Problem 1 suggests that the data were 7 times more likely to be observed under the null hypothesis than the alternative.

To test Hypothesis 2, we conducted a Two-Proportions z -test. We then calculated the effect size, Cohen's d , with a 95% confidence interval (Table 3). The results of Problem 4 ($z = 4.81$, $p < .001$, $d = 0.30$, CI [0.18,0.43]) and Problem 5 ($z = 6.97$, $p < .001$, $d = 0.45$, CI [0.32,0.57]) supported predictions of the original article that more subjects chose the enriched alternative when asked to choose than when asked to reject. However, more subjects chose the enriched alternative when asked to reject than to choose in Problem 7 ($z = -8.04$, $p = 1.00$, $d = -0.52$, CI [-0.64,-0.39]) and Problem 8 ($z = -4.32$, $p = .999$, $d = -0.22$, CI [-0.32,-0.12]), which contradicted the findings in the original article. We found no support for differences between the proportions of subjects choosing the enriched

alternative in the choosing and rejecting decision frame in Problem 1 ($z = 0.56$, $p = .288$, $d = 0.03$, CI $[-0.09, 0.16]$), Problem 2 ($z = -1.37$, $p = .915$, $d = -0.09$, CI $[-0.21, 0.04]$), Problem 3 ($z = -1.58$, $p = .943$, $d = -0.10$, CI $[-0.22, 0.02]$) and Problem 6 ($z = 1.06$, $p = .144$, $d = 0.07$, CI $[-0.06, 0.19]$).

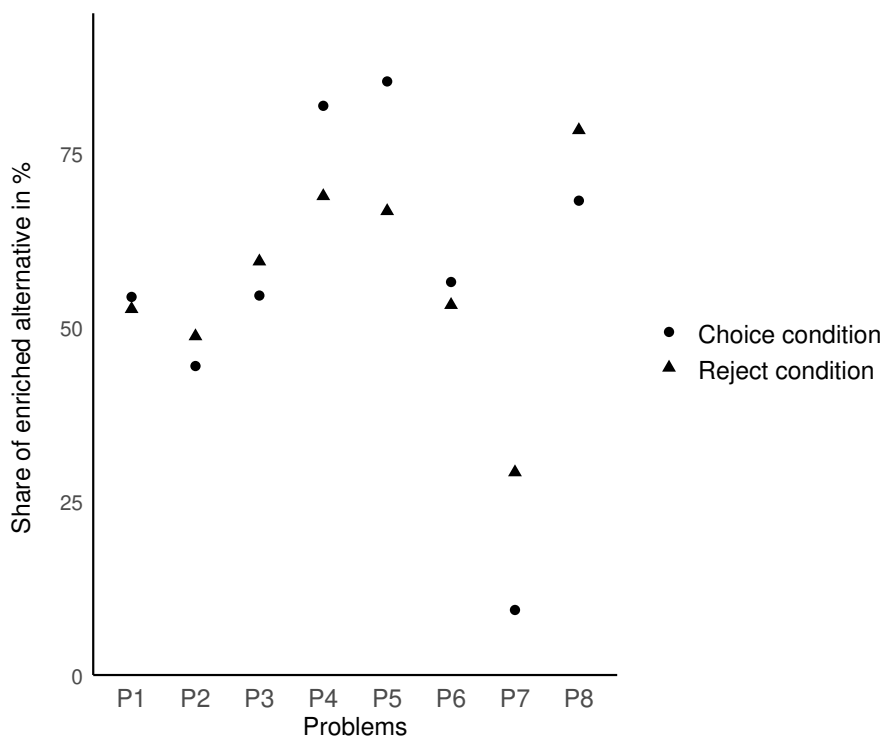


FIGURE 2: Share of the enriched alternative in % between 'choose' and 'reject' experimental conditions.

Furthermore, we conducted Bayesian A/B testing that mirrors the two-proportion z-test based on a model by Kass and Vaidyanathan (1992) using the 'abtest' R package. Mirroring Hypothesis 1, the results for Hypothesis 2 revealed that the Bayes factor (BF) for Problems 1, 2, 3, 6, 7, and 8 are more in favor of the null: Problem 1: $BF_{10} = 0.21$, $BF_{01} = 4.85$; Problem 2: $BF_{10} = 0.05$, $BF_{01} = 18.51$; Problem 3: $BF_{10} = 0.05$, $BF_{01} = 19.89$; Problem 6: $BF_{10} = 0.38$, $BF_{01} = 2.65$; Problem 7: $BF_{10} = 0.02$, $BF_{01} = 64.33$; Problem 8: $BF_{10} = 0.02$, $BF_{01} = 46.16$.

3.1 Comparison of the results with the original findings by Shafir (1993)

The evaluation of the replication results by the pairwise comparisons of each of the eight decision scenarios using LeBel et al.'s (2019) framework are summarized in Table 3. The findings of the present replication are mostly inconsistent with the results of Shafir's original study. Only two of the eight problems (Problem 4 and Problem 5) are supportive of the

compatibility hypothesis. Moreover, two other problems (Problem 7 and Problem 8) showed an effect in the opposite direction. Taken together, the replication findings do not indicate consistent support for the original findings.

3.2 General Summary: Mini meta-analysis

The variations in the findings reported across the eight different decision scenarios make it hard to succinctly summarize the overall effect size of the predictions based on the compatibility hypothesis. Therefore, we conducted a mini meta-analysis of the effect sizes observed across eight decision scenarios for each of the predictions (Goh, Hall & Rosenthal, 2016; Lakens & Etz, 2017). We ran both a within-subject aggregation and a fixed-effects model analysis method using the *'metafor'* package in R (Viechtbauer, 2010; see Figure S3–S6 in the Supplementary Material), and results were near identical.

The mini-meta analysis findings were: Hypothesis 1 $d = -0.01 [-0.06,0.03]$, and Hypothesis 2 $d = -0.01 [-0.06,0.03]$. The results of the mini-meta analysis are summarized in Table 6.

TABLE 6: Summary of findings of the original study versus replication, based on mini-meta analysis.

Predictions	Cohen's d		Replication summary
	Shafir (1993)	Replication	
Hypothesis 1	0.32 [0.23,0.40]	-0.01 [-0.06,0.03]	No signal – inconsistent
Hypothesis 2	0.38 [0.29,0.46]	-0.01 [-0.06,0.03]	No signal – inconsistent

3.3 Extension: Attractiveness ratings

We tested additional variables recorded on a continuous scale that measured the attractiveness of the alternatives. The responses to these additional variables included the attractiveness of the alternatives on a 6-point continuous scale (ranged from 0 to 5). We provide detailed results of the analysis in the Supplementary Materials (see Table S3-S5).

We conducted two sets of independent t-tests. First, we compared the attractiveness of the enriched alternative between the choice and reject experimental conditions. Second, we contrasted the relative attractiveness of the enriched alternatives between the choice and reject experimental conditions. The calculation of the relative attractiveness of the enriched alternative involved subtracting the attractiveness score of the enriched alternative from the attractiveness score of the impoverished alternative within each experimental condition. Then we contrasted the relative attractiveness of the enriched alternative between choice and reject experimental conditions. As Problem 8 included a non-binary alternative, we averaged the attractiveness scores of the impoverished alternatives before calculating the relative attractiveness of the enriched alternative. Furthermore, we conducted Bayesian

analysis for both the planned contrasts with a prior value set at 0.707 (reflecting expectations for an effect, as it was expected from the original study).

The effect size estimates for the enriched alternative's attractiveness in comparing between the choice and reject experimental conditions ranged from 0.00 [−0.12,0.13] to 0.09 [−0.03,0.21]. Furthermore, effect size estimates of relative attractiveness of the enriched alternative across conditions ranged from 0.01 [−0.11,0.14] to 0.14 [0.02,0.26]. The Bayesian analysis mirrors these effect sizes and indicates support for the null in all the problems except Problem 8.

3.4 Extension: Individual-level predictors

We tested the prediction that individuals who rate themselves higher on ability to choose and indicated higher preference for choice are more likely to prefer the enriched alternative. We conducted two separate binary logistic mixed-effects regression analyses which included the experimental condition and individual-level variables as the fixed effect predictors of choosing the enriched alternative (*Yes* = 1; *No* = 0). The regression included subject ID as a random factor on the intercept.

We found no evidence for an association between subjects' ability to choose (Wald $\chi^2(1) = 0.90$, $p = .343$) or subjects' preference for choice (Wald $\chi^2(1) = 0.41$, $p = .522$) and the likelihood of preferring the enriched alternative (see Table S6-S9 for detailed results).

3.5 Extension: Testing the accentuation hypothesis

The inconsistent results regarding the compatibility hypothesis may have been due to the variation of the overall attractiveness of the enriched alternative relative to the impoverished alternative across the eight problems. The accentuation hypothesis (Wedell, 1997) proposed that if the overall relative attractiveness of the enriched alternative is greater than that of the impoverished alternative in a choice set, the positive attributes are more accentuated in the choice condition compared to the reject condition, because of a greater need for justification in the choice condition. Therefore, people more often prefer the enriched alternative in the choice condition than in the reject condition. In contrast, when the overall relative attractiveness of the enriched alternative is lower than that of the impoverished alternative, the negative attributes are more accentuated in the choice condition, again due to greater need for justification. Therefore, in this scenario, people prefer the impoverished alternative in the choice condition more often than in the reject condition.

To test the accentuation hypothesis, we conducted binary logistic mixed-effects regression analysis. In this analysis, we included responses from Problem 1 to 7, as these problems shared the common procedure of choosing between two alternatives (binary choice set). We followed Wedell's (1997) approach to calculate the overall proportion of subjects (across experimental conditions) preferring the enriched alternative for each of the seven problems, as a measure of the overall relative attractiveness of the enriched alternative. We conducted

a binary logistic mixed-effects regression analysis in which the experimental condition, the overall proportion preferring the enriched alternative, and the interaction term (overall proportions \times experimental condition) were the fixed effects predictors of choosing the enriched alternative (*Yes* = 1; *No* = 0). The regression included subject ID as a random factor on the intercept.

The results of the regression found the main effect of the overall proportion preferring the enriched alternative as Wald $\chi^2(1) = 657.28$, $p < .001$, and the interaction effect Wald $\chi^2(1) = 127.70$, $p < .001$ (also see Table 7). As can be seen in Figure 3, the proportions preferring the enriched alternative for the choice and reject experimental conditions as a function of the overall proportion preferring the enriched alternative indicate alternate paths. Across 7 problems, the overall proportion preferring the enriched alternative ranged from 19% to 76%, and the results are consistent with the accentuation hypothesis.

TABLE 7: Results of binary logistic mixed-effects regression following Wedell's (1997) procedure.

	Dependent variable: Predicted probability of enriched alternative	
	Main effect	Interaction
Constant	-2.34*** (0.100)	-3.59*** (0.163)
Overall proportion preferring enriched (PEN)	4.69*** (0.168)	6.97*** (0.286)
Experimental condition (EXP) (1 = Choose; 0 = Reject)	-0.05 (0.058)	2.12*** (0.201)
PEN \times EXP		-3.95*** (0.350)
Observations	7,182	7,182
Log Likelihood	-4,455.54	-4,387.08
Akaike Inf. Crit.	8,919.07	8,784.15
Bayesian Inf. Crit.	8,946.59	8,818.55

Note: * $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$

We also tested the accentuation hypothesis using the attractiveness measures. For each subject we calculated the relative attractiveness of the enriched alternative by subtracting the attractiveness score of the enriched alternative from the attractiveness score of the impoverished alternative across the seven binary problems. This analysis allowed us to test the accentuation hypothesis in a fine-grained manner by taking into account the relative attractiveness measure at the subject level for each of the seven decision problems.

We then conducted a binary logistic mixed-effects regression analysis in which the experimental condition, the relative attractiveness of the enriched alternative, and the interaction term (relative attractiveness of the enriched alternative \times experimental condition) were the fixed effects predictors of choosing the enriched alternative (*Yes* = 1; *No* = 0). The regression included subject ID as a random factor on the intercept.

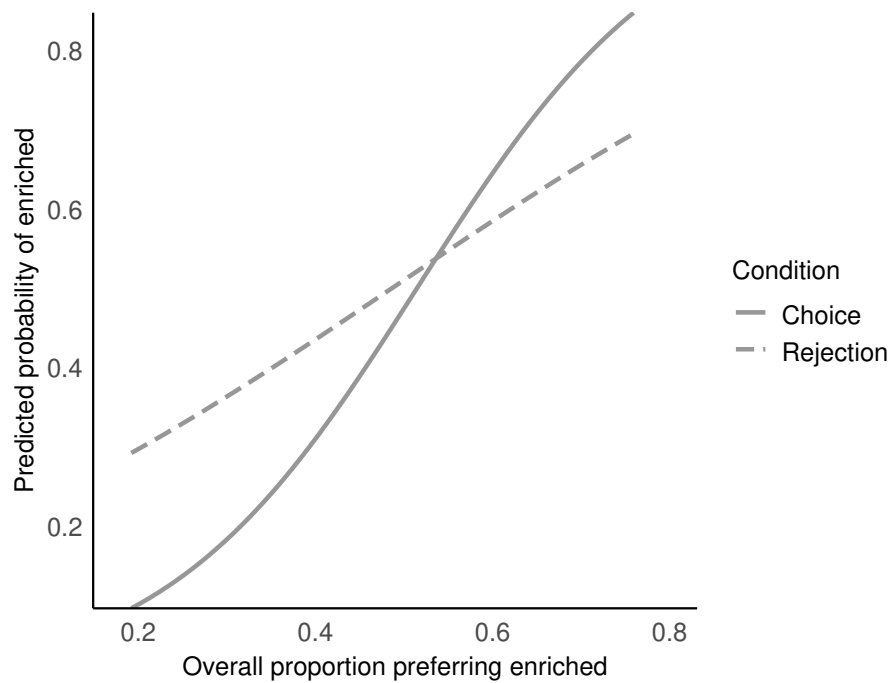


FIGURE 3: Predicted probability of the enriched alternative in choice and rejection tasks as a function of overall preference for the enriched alternative. Fitted lines are the marginal effects of interaction terms.

TABLE 8: Results of binary logistic mixed-effects regression.

	Dependent variable: Predicted probability of enriched alternative	
	Main effect	Interaction
Constant	0.35*** (0.039)	0.41*** (0.043)
Relative attractiveness of enriched alternative (AEO)	0.71*** (0.021)	1.01*** (0.037)
Experimental condition (EXP) (1 = Choose; 0 = Reject)	-0.12** (0.055)	-0.20*** (0.057)
AEO × EXP		-0.51*** (0.044)
Observations	7,182	7,182
Log Likelihood	-4,115.19	-4,043.28
Akaike Inf. Crit.	8,238.37	8,096.55
Bayesian Inf. Crit.	8,265.89	8,130.95

Note: * $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$. The relative attractiveness variable used in the regression was calculated based on the responses to extension variables.

This analysis showed a main effect of the relative attractiveness of enriched alternative (Wald $\chi^2(1) = 980.04, p < .001$) and the interaction term (Wald $\chi^2(1) = 134.08, p < .001$). As can be seen in Figure 4 (also see Table 8), the proportions preferring the enriched alternative for choice and reject experimental conditions as a function of the relative attractiveness of the enriched alternative indicate alternating paths. In summary, the results are consistent with the accentuation hypothesis.

Furthermore, we conducted additional analysis to check the robustness of the results by accounting for the sampling variability of the stimuli (Judd, Westfall & Kenny, 2012). We conducted the same two sets of mixed-effect regression analyses with additional random intercepts and random condition slopes for stimuli along with other predictors. The results of the additional analysis remain the same (see Table S10–S11).

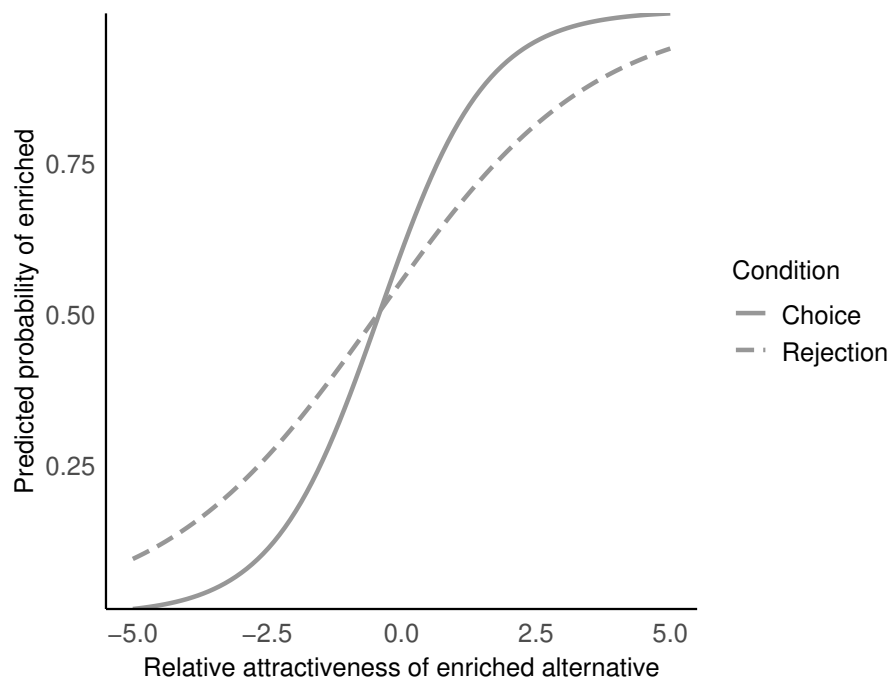


FIGURE 4: Predicted probability of the enriched alternative in choice and rejection tasks as a function of overall preference for the enriched alternative. Fitted lines are the marginal effects of interaction terms. The relative attractiveness variable used in the regression was calculated based on the responses to extension variables.

4 Discussion

We conducted a replication of the eight choosing versus rejection problems in Shafir (1993). We successfully replicated the results of Problem 4 and Problem 5 of the original study. However, in Problems 7 and 8 we found effects in the direction opposite to the original findings and our findings for Problems 1, 2, 3, and 6 indicated support for the null hypothesis. Taken together, we failed to find consistent support for the compatibility hypothesis noted

in Shafir (1993). Additionally, we conducted supplementary analyses and found support for the accentuation hypothesis.

4.1 Replications: Adjustments, implications, and future directions

We aimed for a very close direct replication of the original study, with minimum adjustments, addressing many of the concerns raised over the replication by Many Labs 2, yet our replication still differed from the original studies in several ways. The stimuli used in the original article were targeted at and tested with American undergraduates in the context of the 1990s. We ran the same materials, with no adjustments to the stimuli, online, and with a more diverse population. We also made adjustments to the procedures, by presenting our subjects with all eight questions, instead of only two or three as in the original study, and with no filler items (see Table 4). Replications are never perfectly exact, and given these changes, it is possible that these factors may have somehow affected the results. However, our findings were not random, but rather demonstrated a pattern of results that replicated findings from a different article and supporting an alternative account, and we believe it is highly unlikely that such a change could be explained by any of the adjustments we made.

We also note limitations that suggest promising directions for future research. First, there is a limitation in the extension analysis in that the use of attractiveness rating in testing compatibility hypothesis is not theoretically precise for testing the predictions of compatibility hypothesis. We would also like to see further work testing nuanced preferences (rather than binary choice) yet with more explicit direct integration with the choose/reject framing. Second, it is possible that the inconsistent findings regarding the compatibility hypothesis are due to deviations of auxiliary theories embedded in the compatibility hypothesis (Meehl, 1990). For example, the compatibility hypothesis spells out the substantive argument that people seek positive reasons to justify choosing an alternative, and negative reasons to justify rejecting an alternative. However, auxiliary theories that specify the degree to which justification is a component of the compatibility hypothesis are not well specified and are not clear (as unfortunately is standard in our field). We call for researchers to specify more precise indicators of the boundary conditions of theory testing, so that if some of the contextual factors change, we would be able to directly test and analyze how these affect our findings, rather than engage in post hoc theorizing. Thus, our findings may be due to changes in the conjunction of several premises assumed around the compatibility hypothesis' substantive theory, yet we need stronger well defined theories and hypotheses, and continuous testing over time, to be able to truly assess if and to what extent any of these factors are indeed relevant to the theory, and to the empirical test that theory.

The current study contributes to the theory development by qualifying the theoretical assertions of the compatibility hypothesis. We addressed the methodological issues raised by Shafir (2018) in his commentary on the Many Labs 2 replication. Given our findings, we believe that most explanations noted in the commentary are unlikely reasons for the failure to replicate reported in Many Labs 2 or our failure to find consistent support for the original

findings and the compatibility hypothesis. Theoretical accounts need well-defined criteria that would allow for falsification of these accounts, and our replications help advance theory by testing theoretical assertions of the compatibility hypothesis (Popper, 2002). By improving on the design of Many Labs 2, and by conducting extensions that showed support for plausible alternative accounts, our replication contributes to theory specification and supports further theory development (Glöckner & Betsch, 2011). Researchers conducting research in this domain and future research on this phenomenon can build on insights gained here to advance theory by defining the boundary conditions under which it operates and explore further ways on how it should be tested. Our replication does not rule out the compatibility account, only indicates that it is in need of further elaboration and specification, and further testing, and we see much promise in examining the interaction of the two accounts.

We tested the competing theoretical assertion by Wedell (1997). Our results in support of this account suggest that the stimuli from the 1990s are still of relevance, at least for testing that account. It is still possible that other stimuli developed using the choosing versus rejecting paradigm may show support for the compatibility hypothesis reported by Shafir (1993). Yet, given the Many Labs 2 and our findings we recommend that other compatibility hypothesis stimuli be revisited with direct close replications or that new stimuli be developed before further expanding on the compatibility hypothesis. For this phenomenon, and the judgement and decision-making literature overall, we see great value in conducting well-powered, preregistered direct replications, preferably in Registered Reports or blinded outcomes peer review format. Our findings suggest that future work on choosing versus rejecting may benefit from paying closer attention to the accentuation hypothesis (Wedell, 1997).

4.2 Importance of direct replications

This replication case study highlights the importance of conducting comprehensive direct replications. Many Labs 2 was one of the largest replication efforts to date, yet such mass collaboration replication efforts cannot and should not be taken as a replacement for singular comprehensive direct replications. These large replication projects are valuable in targeting specific research questions about the overall replicability of a research domain, and investigating factors such as heterogeneity and high-level moderators such as culture or setting. Furthermore, large replication projects tend to summarize complex replications in simplified conclusions that fail to capture the complexity inherent in the original articles or the richness of the original and the replication's findings. Therefore, we believe that large scale replication projects should be complemented by singular direct replication and extension studies such as the one we conducted here. Combined, they can help better understand the phenomenon of interest and inform future research.

References

- Brandt, M. J., IJzerman, H., Dijksterhuis, A., Farach, F. J., Geller, J., Giner-Sorolla, R., ... & Van't Veer, A. (2014). The replication recipe: What makes for a convincing replication?. *Journal of Experimental Social Psychology, 50*, 217–224. 10.1016/j.jesp.2013.10.005
- Chernev, A. (2009). Choosing versus rejecting: The impact of goal-task compatibility on decision confidence. *Social Cognition, 27*(2), 249–260. <https://doi.org/10.1521/soco.2009.27.2.249>
- Coles, N. A., Tiokhin, L., Scheel, A. M., Isager, P. M., & Lakens, D. (2018). The costs and benefits of replication studies. *Behavioral and Brain Sciences, 41*, e124. <https://doi.org/10.1017/S0140525X18000596>
- Feldman, G., Baumeister, R. F., & Wong, K. F. E. (2014). Free will is about choosing: The link between choice and the belief in free will. *Journal of Experimental Social Psychology, 55*, 239–245. 10.1016/j.jesp.2014.07.012
- Ganzach, Y. (1995). Attribute scatter and decision outcome Judgment versus choice. *Organizational Behavior & Human Decision Processes, 62*(1), 113–122. 10.1006/obhd.1995.1036
- Glöckner, A., & Betsch, T. (2011). The empirical content of theories in judgment and decision making: Shortcomings and remedies. *Judgment and Decision Making, 6*(8), 711–721.
- Goh, J. X., Hall, J. A., & Rosenthal, R. (2016). Mini meta-analysis of your own studies: Some arguments on why and a primer on how. *Social and Personality Psychology Compass, 10*(10), 535–549. 10.1111/spc3.12267
- Isager, P. M. (2019). Quantifying replication value: A formula-based approach to study selection in replication research. *Open Science 2019*, Trier, Germany. 10.23668/psycharchives.2392. See "Quantifying the corroboration of a finding" preprint retrieved from: <https://pedermisager.netlify.com/post/quantifying-the-corroboration-of-a-finding/>
- Judd, C. M., Westfall, J., & Kenny, D. A. (2012). Treating stimuli as a random factor in social psychology: A new and comprehensive solution to a pervasive but largely ignored problem. *Journal of personality and social psychology, 103*(1), 54–69. 10.1037/a0028347
- Kahneman, D. (2003). A perspective on judgment and choice: mapping bounded rationality. *American psychologist, 58*(9), 697–720. 10.1037/0003--066X.58.9.697
- Kass, R. E., & Vaidyanathan, S. K. (1992). Approximate Bayes factors and orthogonal parameters, with application to testing equality of two binomial proportions. *Journal of the Royal Statistical Society: Series B (Methodological), 54*(1), 129–144. 10.1111/j.2517--6161.1992.tb01868.x
- Klein, R. A., Vianello, M., Hasselman, F., Adams, B. G., Adams, R. B., Jr., Alper, S., . . . & Nosek, B. A. (2018). Many Labs 2: Investigating variation in replicability across samples and settings. *Advances in Methods and Practices in Psychological Science, 1*(4), 443–490. 10.1177/2515245918810225

- Kruschke, J. K., & Liddell, T. M. (2018). The Bayesian new statistics: Hypothesis testing, estimation, meta-analysis, and power analysis from a Bayesian perspective. *Psychonomic Bulletin & Review*, *25*, 178–206. 10.3758/s13423--016--1221--4
- Lakens, D., & Etz, A. J. (2017). Too true to be bad: When sets of studies with significant and nonsignificant findings are probably true. *Social Psychological and Personality Science*, *8*(8), 875–881. 10.1177/1948550617693058
- LeBel, E. P., McCarthy, R. J., Earp, B. D., Elson, M., & Vanpaemel, W. (2018). A unified framework to quantify the credibility of scientific findings. *Advances in Methods and Practices in Psychological Science*, *1*(3), 389–402. 10.1177/2515245918787489
- LeBel, E. P., Vanpaemel, W., Cheung, I., & Campbell, L. (2019). A brief guide to evaluate replications. *Meta Psychology*, *3*, 1–17. 10.15626/MP.2018.843
- Litman, L., Robinson, J., & Abberbock, T. (2017). TurkPrime. com: A versatile crowdsourcing data acquisition platform for the behavioral sciences. *Behavior Research Methods*, *49*, 433–442. 10.3758/s13428--016--0727-z
- Meehl, P. E. (1990). Appraising and amending theories: The strategy of Lakatosian defense and two principles that warrant it. *Psychological inquiry*, *1*(2), 108–141.
- Morewedge, C. K., & Kahneman, D. (2010). Associative processes in intuitive judgment. *Trends in cognitive sciences*, *14*(10), 435–440. DOI: 10.1016/j.tics.2010.07.004
- Morey, R. D., & Rouder, J. N. (2015). BayesFactor: Computation of Bayes factors for common designs (R Package Version 0.9.12–2). Retrieved from <https://CRAN.R-project.org/package=BayesFactor>
- Nagpal, A., & Krishnamurthy, P. (2008). Attribute conflict in consumer decision making: The role of task compatibility. *Journal of Consumer Research*, *34* (5), 696–705. <http://doi.org/10.1086/521903>
- Open Science Collaboration (2015). Estimating the reproducibility of psychological science. *Science*, *349* (6251), aac4716. 10.1126/science.aac4716
- Park, C. W., Jun, S. Y., & MacInnis, D. J. (2000). Choosing what I want versus rejecting what I do not want: An application of decision framing to product option choice decisions. *Journal of Marketing Research*, *37*(2), 187–202. 10.1509/jmkr.37.2.187.18731
- Popper, K. (2002). *The logic of scientific discovery (2nd ed.)*. London, England: Routledge.
- R Core Team (2015). R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria. <http://www.R-project.org>
- Savani, K., Markus, H. R., Naidu, N. V. R., Kumar, S., & Berlia, N. (2010). What counts as a choice? US Americans are more likely than Indians to construe actions as choices. *Psychological Science*, *21*(3), 391–398. DOI: 10.1177/0956797609359908
- Shafir, E. (1993). Choosing versus rejecting: Why some options are both better and worse than others. *Memory & cognition*, *21*(4), 546–556.
- Shafir, E. (2018). The workings of choosing and rejecting: Commentary on Many Labs 2. *Advances in Methods and Practices in Psychological Science*, *1*(4), 495–496. 10.1177/2515245918814812

- Slovic, P., Griffin, D., & Tversky, A. (1990). *Compatibility effects in judgment and choice. Insights in decision making: Theory and applications* (pp. 5–27). Chicago, IL: University of Chicago Press.
- Sokolova, T., & Krishna, A. (2016). Take it or leave it: How choosing versus rejecting alternatives affects information processing. *Journal of Consumer Research*, *43*(4), 614–635. 10.1093/jcr/ucw049
- Tversky, A., & Kahneman, D. (1986). Rational choice and the framing of decisions. *Journal of Business*, *59* (4), 251–278.
- Tversky, A., Sattath, S., & Slovic, P. (1988). Contingent weighting in judgment and choice. *Psychological Review*, *95*(3), 371–384. 10.1037/0033--295X.95.3.371
- Vandekerckhove, J., Rouder, J. N., & Kruschke, J. K. (2018). Editorial: Bayesian methods for advancing psychological science. *Psychonomic Bulletin & Review*, *25*, 1–4. 10.3758/s13423--018--1443--8
- van't Veer, A.E., & Giner-Sorolla, R. (2016). Pre-registration in social psychology—A discussion and suggested template. *Journal of Experimental Social Psychology*, *67*, 2–12. 10.1016/j.jesp.2016.03.004
- Viechtbauer, W. (2010). Conducting meta-analyses in R with the metafor package. *Journal of Statistical Software*, *36*(3), 1–48. <http://www.jstatsoft.org/v36/i03/>
- von Neumann, J., & Morgenstern, O. (1947). *Theory of games and economic behavior*. Princeton, NJ: Princeton University Press.
- Wedell, D. H. (1997). Another look at reasons for choosing and rejecting. *Memory & Cognition*, *25*(6), 873–887. 10.3758/BF03211332
- Zwaan, R. A., Etz, A., Lucas, R. E., & Donnellan, M. B. (2018). Making replication mainstream. *Behavioral and Brain Sciences*, *41*, E120. 10.1017/S0140525X17001972